# Projection, matching, and basis pursuits for multilinear approximations

Lek-Heng Lim

University of California, Berkeley

July 10, 2008

# Best *r*-term approximation

$$f \approx \alpha_1 f_1 + \alpha_2 f_2 + \cdots + \alpha_r f_r.$$

- **Target function** $f \in \mathcal{H}$ vector space, cone, etc.
- $f_1, \ldots, f_r \in \mathscr{D} \subset \mathcal{H}$ **dictionary**.
- $\alpha_1, \ldots, \alpha_r \in \mathbb{R}$ or $\mathbb{C}$ (linear), $\mathbb{R}_+$ (convex), $\mathbb{R} \cup \{-\infty\}$ (tropical).
- $\approx$ with respect to $\varphi : \mathcal{H} \times \mathcal{H} \to \mathbb{R}$, some measure of 'nearness' between pairs of points (e.g. norms, metric, volumes, expectation, entropy, Brègman divergences, etc), want

$$\mathrm{argmin}\{\varphi(f, \alpha_1 f_1 + \ldots \alpha_r f_r) \mid f_i \in \mathscr{D}\}.$$

- For concreteness, $\mathcal{H}$ separable Hilbert space; measure of nearness is a norm, but not necessarily the one induced by its inner product.
- Reference: various papers by A. Cohen, R. DeVore, V. Temlyakov.

## Dictionaries

- Number base: $\mathscr{D} = \{10^n \mid n \in \mathbb{Z}\} \subseteq \mathbb{R}$,

$$\tfrac{22}{7} = 3 \cdot 10^0 + 1 \cdot 10^{-1} + 4 \cdot 10^{-2} + 2 \cdot 10^{-3} + \cdots$$

- Spanning set: $\mathscr{D} = \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\} \subseteq \mathbb{R}^2$,

$$\begin{bmatrix} 2 \\ -3 \end{bmatrix} = 3 \begin{bmatrix} 1 \\ -1 \end{bmatrix} - 1 \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

- Taylor: $\mathscr{D} = \{x^n \mid n \in \mathbb{N} \cup \{0\}\} \subseteq C^\omega(\mathbb{R})$,

$$\exp(x) = 1 + x + \tfrac{1}{2}x^2 + \tfrac{1}{6}x^3 + \cdots$$

- Fourier: $\mathscr{D} = \{\cos(nx), \sin(nx) \mid n \in \mathbb{Z}\} \subseteq L^2(-\pi, \pi)$,

$$\tfrac{1}{2}x = \sin(x) - \tfrac{1}{2}\sin(2x) + \tfrac{1}{3}\sin(3x) - \cdots$$

- $\mathscr{D}$ orthonormal basis, Schauder basis, Hamel basis, Riesz basis, frames, a dense spanning set.

## More dictionaries

- Discrete cosine:

$$\mathscr{D} = \left\{ \sqrt{\tfrac{2}{N}} \cos(k + \tfrac{1}{2})(n + \tfrac{1}{2})\tfrac{\pi}{N} \;\middle|\; k \in [N-1] \right\} \subseteq \mathbb{C}^N.$$

- Peter-Weyl:

$$\mathscr{D} = \{\langle \pi(x)\mathbf{e}_i, \mathbf{e}_j \rangle \mid \pi \in \widehat{G}, i, j \in [d_\pi]\} \subseteq L^2(G).$$

- Paley-Wiener:

$$\mathscr{D} = \{\operatorname{sinc}(x - n) \mid n \in \mathbb{Z}\} \subseteq H^2(\mathbb{R}).$$

- Gabor:

$$\mathscr{D} = \{e^{i\alpha n x} e^{-(x - m\beta)^2/2} \mid (m, n) \in \mathbb{Z} \times \mathbb{Z}\} \subseteq L^2(\mathbb{R}).$$

- Wavelet:

$$\mathscr{D} = \{2^{n/2}\psi(2^n x - m) \mid (m, n) \in \mathbb{Z} \times \mathbb{Z}\} \subseteq L^2(\mathbb{R}).$$

- Friends of wavelets: $\mathscr{D} \subseteq L^2(\mathbb{R}^2)$ beamlets, brushlets, curvelets, ridgelets, wedgelets, multiwavelets [Mohlenkamp, Pereyra; 2008].

# Approximants

## Definition

Dictionary $\mathscr{D} \subset \mathcal{H}$. For $r \in \mathbb{N}$, the set of **r-term approximants** is

$$\Sigma_r(\mathscr{D}) := \left\{ \sum_{i=1}^{r} \alpha_i f_i \in \mathcal{H} \;\middle|\; \alpha_i \in \mathbb{C}, f_i \in \mathscr{D} \right\}.$$

Let $f \in \mathcal{H}$. The **error of r-term approximation** is

$$\sigma_n(f) := \inf_{g \in \Sigma_r(\mathscr{D})} \|f - g\|.$$

- Linear combination of two $r$-term approximants may have more than $r$ non-zero terms.
- $\Sigma_r(\mathscr{D})$ not a subspace of $\mathcal{H}$. Hence **nonlinear approximation**.
- In contrast with usual (linear) approximation, ie.

$$\inf_{g \in \mathsf{span}(\mathscr{D})} \|f - g\|.$$

# Small is beautiful

$$f \approx \sum_{i \in \mathscr{I} \subseteq \mathscr{D}} \alpha_i f_i$$

- Want good approximation, ie. $\|f - \sum_{i \in \mathscr{I} \subseteq \mathscr{D}} \alpha_i f_i\|$ small.
- Want sparse/concentrated representation, ie. $|\mathscr{I}|$ small.
- Sparsity depends on choice of $\mathscr{D}$.

   ▶ $\mathscr{D}_{10} = \{10^n \mid n \in \mathbb{Z}\}, \mathscr{D}_3 = \{3^n \mid n \in \mathbb{Z}\} \subseteq \mathbb{R}$,

   $$\begin{aligned} \tfrac{1}{3} &= [0.33333\cdots]_{10} = \sum_{n=1}^{\infty} 3 \cdot 10^{-n} \\ &= [0.1]_3 = 1 \cdot 3^{-1}. \end{aligned}$$

   ▶ $\mathscr{D}_{\text{fourier}} = \{\cos(nx), \sin(nx) \mid n \in \mathbb{Z}\}$,

   $$\tfrac{1}{2}x = \sin(x) - \tfrac{1}{2}\sin(2x) + \tfrac{1}{3}\sin(3x) - \cdots.$$

   ▶ $\mathscr{D}_{\text{taylor}} = \{x^n \mid n \in \mathbb{N} \cup \{0\}\}$,

   $$\sin(x) = x - \tfrac{1}{6}x^3 + \tfrac{1}{120}x^5 - \cdots.$$

# Bigger is better

- **Union of dictionaries:** allows for efficient (sparse) representation of different features
  - $\mathscr{D} = \mathscr{D}_{\text{fourier}} \cup \mathscr{D}_{\text{wavelets}}$,
  - $\mathscr{D} = \mathscr{D}_{\text{spikes}} \cup \mathscr{D}_{\text{sinusoids}} \cup \mathscr{D}_{\text{splines}}$,
  - $\mathscr{D} = \mathscr{D}_{\text{wavelets}} \cup \mathscr{D}_{\text{curvelets}} \cup \mathscr{D}_{\text{beamlets}} \cup \mathscr{D}_{\text{ridgelets}}$.

- $\mathscr{D}$ **overcomplete** or **redundant** dictionary. Trade off: computational complexity.

- **Rule of thumb:** the larger and more diverse the dictionary, the more efficient/sparser the representation.

- **Observation:** $\mathscr{D}$ above all zero dimensional (at most countably infinite).

- **Question:** What about dictionaries with a continuously varying families of functions?

- **Meta question:** Why should tensor folks care about this?

# Vectors, matrices, tensors: functions on finite sets

Totally ordered finite sets: $[n] = \{1 < 2 < \cdots < n\}$, $n \in \mathbb{N}$.

- Vector or $n$-tuple

$$f : [n] \to \mathbb{R}.$$

  If $f(i) = a_i$, then $f$ is represented by $\mathbf{a} = [a_1, \ldots, a_n]^\top \in \mathbb{R}^n$.

- Matrix

$$f : [m] \times [n] \to \mathbb{R}.$$

  If $f(i,j) = a_{ij}$, then $f$ is represented by $A = [a_{ij}]_{i,j=1}^{m,n} \in \mathbb{R}^{m \times n}$.

- Hypermatrix (order 3)

$$f : [l] \times [m] \times [n] \to \mathbb{R}.$$

  If $f(i,j,k) = a_{ijk}$, then $f$ is represented by $\mathcal{A} = [\![a_{ijk}]\!]_{i,j,k=1}^{l,m,n} \in \mathbb{R}^{l \times m \times n}$.

Normally $\mathbb{R}^X = \{f : X \to \mathbb{R}\}$. Ought to be $\mathbb{R}^{[n]}, \mathbb{R}^{[m] \times [n]}, \mathbb{R}^{[l] \times [m] \times [n]}$.

# Hilbert space structure

- $\ell^2([n])$: $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$, $\langle \mathbf{a}, \mathbf{b} \rangle = \mathbf{a}^\top \mathbf{b} = \sum_{i=1}^{n} a_i b_i$.
- $\ell^2([m] \times [n])$: $A, B \in \mathbb{R}^{m \times n}$, $\langle A, B \rangle = \mathrm{tr}(A^\top B) = \sum_{i,j=1}^{m,n} a_{ij} b_{ij}$.
- $\ell^2([l] \times [m] \times [n])$: $\mathcal{A}, \mathcal{B} \in \mathbb{R}^{l \times m \times n}$, $\langle \mathcal{A}, \mathcal{B} \rangle = \sum_{i,j,k=1}^{l,m,n} a_{ijk} b_{ijk}$.
- In general,

$$\ell^2([m] \times [n]) = \ell^2([m]) \otimes \ell^2([n]),$$
$$\ell^2([l] \times [m] \times [n]) = \ell^2([l]) \otimes \ell^2([m]) \otimes \ell^2([n]).$$

- Frobenius norm

$$\|\mathcal{A}\|_F^2 = \sum_{i,j,k=1}^{l,m,n} a_{ijk}^2.$$

## Hypermatrices and tensors

Up to choice of bases

- $\mathbf{a} \in \mathbb{C}^n$ can represent a vector in $V$ (contravariant) or a linear functional in $V^*$ (covariant).
- $A \in \mathbb{C}^{m \times n}$ can represent a bilinear form $V \times W \to \mathbb{C}$ (contravariant), a bilinear form $V^* \times W^* \to \mathbb{C}$ (covariant), or a linear operator $V \to W$ (mixed).
- $\mathcal{A} \in \mathbb{C}^{l \times m \times n}$ can represent trilinear form $U \times V \times W \to \mathbb{C}$ (contravariant), bilinear operators $V \times W \to U$ (mixed), etc.

A hypermatrix is the same as a tensor if

1. we give it coordinates (represent with respect to some bases);
2. we ignore covariance and contravariance.

## Tensor ranks

- For $\mathbf{u} \in \mathbb{R}^l, \mathbf{v} \in \mathbb{R}^m, \mathbf{w} \in \mathbb{R}^n$,

$$\mathbf{u} \otimes \mathbf{v} \otimes \mathbf{w} := [\![ u_i v_j w_k ]\!]_{i,j,k=1}^{l,m,n} \in \mathbb{R}^{l \times m \times n}.$$

- **Outer product rank.** $\mathcal{A} \in \mathbb{R}^{l \times m \times n}$,

$$\text{rank}_{\otimes}(\mathcal{A}) = \min\{ r \mid \mathcal{A} = \sum_{i=1}^{r} \sigma_i \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i, \quad \sigma_i \in \mathbb{R} \}.$$

- **Symmetric outer product rank.** $\mathcal{A} \in S^k(\mathbb{R}^n)$,

$$\text{rank}_S(\mathcal{A}) = \min\{ r \mid \mathcal{A} = \sum_{i=1}^{r} \lambda_i \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{v}_i, \quad \lambda_i \in \mathbb{R} \}.$$

- **Nonnegative outer product rank.** $\mathcal{A} \in \mathbb{R}_+^{l \times m \times n}$,

$$\text{rank}_+(\mathcal{A}) = \min\{ r \mid \mathcal{A} = \sum_{i=1}^{r} \delta_i \mathbf{x}_i \otimes \mathbf{y}_i \otimes \mathbf{z}_i, \quad \delta_i \in \mathbb{R}_+ \}.$$

# SVD, EVD, NMF of a matrix

- **Singular value decomposition** of $A \in \mathbb{R}^{m \times n}$,

$$A = U\Sigma V^\top = \sum_{i=1}^{r} \sigma_i \mathbf{u}_i \otimes \mathbf{v}_i$$

where $\text{rank}(\mathcal{A}) = r$, $U \in \mathrm{O}(m)$ left singular vectors, $V \in \mathrm{O}(n)$ right singular vectors, $\Sigma$ singular values.

- **Symmetric eigenvalue decomposition** of $A \in \mathrm{S}^2(\mathbb{R}^n)$,

$$A = V\Lambda V^\top = \sum_{i=1}^{r} \lambda_i \mathbf{v}_i \otimes \mathbf{v}_i,$$

where $\text{rank}(A) = r$, $V \in \mathrm{O}(n)$ eigenvectors, $\Lambda$ eigenvalues.

- **Nonnegative matrix factorization** of $A \in \mathbb{R}_+^{n \times n}$,

$$A = X\Delta Y^\top = \sum_{i=1}^{r} \delta_i \mathbf{x}_i \otimes \mathbf{y}_i$$

where $\text{rank}_+(A) = r$, $X, Y \in \mathbb{R}_+^{m \times r}$ unit column vectors (in the 1-norm), $\Delta$ positive values.

# SVD, EVD, NMF of a hypermatrix

- **Outer product decomposition** of $\mathcal{A} \in \mathbb{R}^{l \times m \times n}$,

$$\mathcal{A} = \sum_{i=1}^{r} \sigma_i \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i$$

  where $\mathrm{rank}_{\otimes}(\mathcal{A}) = r$, $\mathbf{u}_i \in \mathbb{R}^l, \mathbf{v}_i \in \mathbb{R}^m, \mathbf{w}_i \in \mathbb{R}^n$ unit vectors, $\sigma_i \in \mathbb{R}$.

- **Symmetric outer product decomposition** of $\mathcal{A} \in S^3(\mathbb{R}^n)$,

$$\mathcal{A} = \sum_{i=1}^{r} \lambda_i \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{v}_i$$

  where $\mathrm{rank}_S(A) = r$, $\mathbf{v}_i$ unit vector, $\lambda_i \in \mathbb{R}$.

- **Nonnegative outer product decomposition** for hypermatrix $\mathcal{A} \in \mathbb{R}_+^{l \times m \times n}$ is

$$\mathcal{A} = \sum_{i=1}^{r} \delta_i \mathbf{x}_i \otimes \mathbf{y}_i \otimes \mathbf{z}_i$$

  where $\mathrm{rank}_+(A) = r$, $\mathbf{x}_i \in \mathbb{R}_+^l, \mathbf{y}_i \in \mathbb{R}_+^m, \mathbf{z}_i \in \mathbb{R}_+^n$ unit vectors, $\delta_i \in \mathbb{R}_+$.

# Best low rank approximation of a matrix

- Given $A \in \mathbb{R}^{m \times n}$. Want

$$\mathrm{argmin}_{\mathrm{rank}(B) \leq r} \|A - B\|.$$

- More precisely, find $\sigma_i, \mathbf{u}_i, \mathbf{v}_i, \ i = 1, \ldots, r$, that minimizes

$$\|\mathcal{A} - \sigma_1 \mathbf{u}_1 \otimes \mathbf{v}_1 - \sigma_2 \mathbf{u}_2 \otimes \mathbf{v}_2 - \cdots - \sigma_r \mathbf{u}_r \otimes \mathbf{v}_r\|.$$

## Theorem (Eckart–Young)

Let $A = U\Sigma V^\top = \sum_{i=1}^{\mathrm{rank}(A)} \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$ be singular value decomposition. For $r \leq \mathrm{rank}(A)$, let

$$A_r := \sum_{i=1}^{r} \sigma_i \mathbf{u}_i \mathbf{v}_i^\top.$$

Then

$$\|A - A_r\|_F = \min_{\mathrm{rank}(B) \leq r} \|A - B\|_F.$$

- No such thing for hypermatrices of order 3 or higher.

# Segre variety and its secant varieties

- The set of all rank-1 hypermatrices is known as the Segre variety in algebraic geometry.

- It is a closed set (in both the Euclidean and Zariski sense) as it can be described algebraically:

$$\text{Seg}(\mathbb{R}^l, \mathbb{R}^m, \mathbb{R}^n) = \{\mathcal{A} \in \mathbb{R}^{l \times m \times n} \mid \mathcal{A} = \mathbf{u} \otimes \mathbf{v} \otimes \mathbf{w}\} =$$
$$\{\mathcal{A} \in \mathbb{R}^{l \times m \times n} \mid a_{i_1 i_2 i_3} a_{j_1 j_2 j_3} = a_{k_1 k_2 k_3} a_{l_1 l_2 l_3}, \{i_\alpha, j_\alpha\} = \{k_\alpha, l_\alpha\}\}$$

- Hypermatrices that have rank $> 1$ are elements on the higher secant varieties of $\mathscr{S} = \text{Seg}(\mathbb{R}^l, \mathbb{R}^m, \mathbb{R}^n)$.

- E.g. a hypermatrix has rank 2 if it sits on a secant line through two points in $\mathscr{S}$ but not on $\mathscr{S}$, rank 3 if it sits on a secant plane through three points in $\mathscr{S}$ but not on any secant lines, etc.

- Minor technicality: should really be secant *quasiprojective variety*.

## Same thing different names

- $r$th secant (quasiprojective) variety of the Segre variety is the set of $r$ term approximants.
- If $\mathscr{D} = \mathrm{Seg}(\mathbb{R}^l, \mathbb{R}^m, \mathbb{R}^n)$, then

$$\Sigma_r(\mathscr{D}) = \{\mathcal{A} \in \mathbb{R}^{l \times m \times n} \mid \mathrm{rank}_\otimes(\mathcal{A}) \leq r\}.$$

- Rank revealing matrix decompositions (non-unique: LU, QR, SVD):

$$\mathscr{D} = \{\mathbf{x}\mathbf{y}^\top \mid (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^m \times \mathbb{R}^n\} = \{A \in \mathbb{R}^{m \times n} \mid \mathrm{rank}(A) \leq 1\}.$$

- Often unique for tensors [Kruskal; 1977], [Sidiroupoulos, Bro; 2000]:
  - spark$(\mathbf{x}_1, \ldots, \mathbf{x}_r)$ = size of minimal linearly dependent subset [Donoho, Elad; 2003].
  - Decomposition $\mathcal{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i$ is unique up to scaling if

  $$\mathrm{spark}(\mathbf{u}_1, \ldots, \mathbf{u}_r) + \mathrm{spark}(\mathbf{v}_1, \ldots, \mathbf{v}_r) + \mathrm{spark}(\mathbf{w}_1, \ldots, \mathbf{w}_r) \geq 2r + 5.$$

## Dictionaries of positive dimensions

- Neural networks:

$$\mathscr{D} = \{\sigma(\mathbf{w}^\top \mathbf{x} + w_0) \mid (w_0, \mathbf{w}) \in \mathbb{R} \times \mathbb{R}^n\}$$

  where $\sigma : \mathbb{R} \to \mathbb{R}$ sigmoid function, eg. $\sigma(x) = [1 + \exp(-x)]^{-1}$.

- Exponential [Beylkin, Monzón; 2005]:

$$\mathscr{D} = \{e^{-tx} \mid t \in \mathbb{R}_+\} \qquad \text{or} \qquad \mathscr{D} = \{e^{\tau x} \mid \tau \in \mathbb{C}\}.$$

- Outer product decomposition:

$$\mathscr{D} = \{\mathbf{u} \otimes \mathbf{v} \otimes \mathbf{w} \mid (\mathbf{u}, \mathbf{v}, \mathbf{w}) \in \mathbb{R}^l \times \mathbb{R}^m \times \mathbb{R}^n\}$$
$$= \{\mathcal{A} \in \mathbb{R}^{l \times m \times n} \mid \mathrm{rank}_\otimes(\mathcal{A}) \leq 1\}.$$

- Symmetric outer product decomposition:

$$\mathscr{D} = \{\mathbf{v} \otimes \mathbf{v} \otimes \mathbf{v} \mid \mathbf{v} \in \mathbb{R}^n\} = \{\mathcal{A} \in \mathsf{S}^3(\mathbb{R}^n) \mid \mathrm{rank}_\mathsf{S}(\mathcal{A}) \leq 1\}.$$

- Nonnegative outer product decomposition:

$$\mathscr{D} = \{\mathbf{x} \otimes \mathbf{y} \otimes \mathbf{z} \mid (\mathbf{x}, \mathbf{y}, \mathbf{z}) \in \mathbb{R}_+^l \times \mathbb{R}_+^m \times \mathbb{R}_+^n\}$$
$$= \{\mathcal{A} \in \mathbb{R}_+^{l \times m \times n} \mid \mathrm{rank}_+(\mathcal{A}) \leq 1\}.$$

## Pursuit algorithms

- Stepwise projection:

$$g_k = \operatorname{argmin}_{g \in \mathscr{D}}\{\|f - h\| \mid h \in \operatorname{span}\{g_1, \ldots, g_{k-1}, g\}\},$$
$$f_k = \operatorname{proj}_{\operatorname{span}\{g_1, \ldots, g_k\}}(f).$$

- Orthonormal matching pursuit:

$$g_k = \operatorname{argmax}_{g \in \mathscr{D}}|\langle f - f_{k-1}, g \rangle|,$$
$$f_k = \operatorname{proj}_{\operatorname{span}\{g_1, \ldots, g_k\}}(f).$$

- Pure greedy:

$$g_k = \operatorname{argmax}_{g \in \mathscr{D}}|\langle f - f_{k-1}, g \rangle|,$$
$$f_k = f_{k-1} + \langle f - f_{k-1}, g_k \rangle g_k.$$

- Relaxed greedy:

$$g_k = \operatorname{argmin}_{g \in \mathscr{D}}\{\|f - h\| \mid h \in \operatorname{span}\{f_{k-1}, g\}\},$$
$$f_k = \alpha_k f_{k-1} + \beta_k g_k.$$

# Pursuit algorithms for tensor approximations

- Target function

$$f : [l] \times [m] \times [n] \to \mathbb{R}.$$

- Dictionary of **separable functions**,

$$\mathscr{D} = \{g : [l] \times [m] \times [n] \to \mathbb{R} \mid g(i,j,k) = \vartheta(i)\varphi(j)\psi(k)\},$$

where $\vartheta : [l] \to \mathbb{R}$, $\varphi : [m] \to \mathbb{R}$, $\psi : [n] \to \mathbb{R}$.

- Inner product

$$\langle f, g \rangle = \sum_{i,j,k=1}^{l,m,n} f(i,j,k)g(i,j,k).$$

and corresponding norm and projection.

- Ditto for the symmetric and nonnegative versions.
- Details: 11:30am–12:30pm, July 15, 2008, MSRI, Berkeley, CA.

## Advertisement

**Geometry and representation theory of tensors for computer science, statistics, and other areas**

1. MSRI Summer Graduate Workshop

   - July 7 to July 18, 2008
   - Organized by J.M. Landsberg, L.-H. Lim, J. Morton
   - Mathematical Sciences Research Institute, Berkeley, CA
   - `http://msri.org/calendar/sgw/WorkshopInfo/451/show_sgw`

2. AIM Workshop

   - July 21 to July 25, 2008
   - Organized by J.M. Landsberg, L.-H. Lim, J. Morton, J. Weyman
   - American Institute of Mathematics, Palo Alto, CA
   - `http://aimath.org/ARCC/workshops/repnsoftensors.html`